

Extensions of the Markov chain marginal bootstrap[☆]

Masha Kocherginsky^a, Xuming He^{b,*}

^a*Department of Health Studies, University of Chicago, USA*

^b*Department of Statistics, University of Illinois, Urbana-Champaign, USA*

Available online 16 March 2007

Abstract

The Markov chain marginal bootstrap (MCMB) was introduced by He and Hu [2002. Markov chain marginal bootstrap. *J. Amer. Statist. Assoc.* 97(459) (2002) 783–795] as a bootstrap-based method for constructing confidence intervals or regions for a wide class of M -estimators in linear regression and maximum likelihood estimators in certain parametric models. In this article we discuss more general applications of MCMB-A, an extension of the MCMB algorithm, which was first proposed in Kocherginsky et al. [2005. Practical confidence intervals for regression quantiles. *J. Comput. Graphical Statist.* 14, 41–55] for quantile regression models. We also present a further extension of the MCMB algorithm, the **B**-transformation, which is a transformation of the estimating equations, aiming to broaden the applicability of the MCMB algorithm to general estimating equations that are not necessarily likelihood-based. We show that applying the **A**- and **B**-transformations jointly enables the MCMB algorithm to be used for inference related to a very general class of estimating equations. We illustrate the use of the MCMB-AB algorithm with a nonlinear regression model with heteroscedastic error distribution.

© 2007 Elsevier B.V. All rights reserved.

Keywords: Bootstrap; General estimating equations; Confidence intervals; Transformation; Nonlinear regression

1. Introduction

The key idea of the bootstrap is to resample from the empirical distribution of the observed data and subsequently to estimate the quantities of interest (e.g., standard errors and confidence regions) from the bootstrap samples. The bootstrap methods are particularly useful when the asymptotic approximation is hard to compute as it foregoes “the long-winded and error-prone analytical calculation” (Davison and Hinkley, 1997). The bootstrap can, however, be computationally expensive and even prohibitive for large-dimensional problems.

Regression models are among the most important and frequent applications of the bootstrap. Consider a regression model

$$y_i = f(\mathbf{x}_i, \boldsymbol{\theta}) + e_i, \quad (1)$$

[☆]The research is partially supported by the National Science Foundation Award DMS-0604229 and a National Security Agency Grant.

*Corresponding author.

E-mail addresses: mkocherg@uchicago.edu (M. Kocherginsky), x-he@uiuc.edu (X. He).

where the response variable Y is related to the independent variable \mathbf{X} (which can be a matrix) via a function f and a vector of unknown parameters $\boldsymbol{\theta}$. The residual and the paired bootstrap are two common methods for bootstrap estimation of $\boldsymbol{\theta}$ (Efron and Tibshirani, 1998). The residual bootstrap works well when the errors are exchangeable in distribution, while the paired bootstrap is based on the assumption that the observations come from a multivariate distribution $F(\mathbf{X}, \mathbf{Y})$. Either residuals or pairs of observations (\mathbf{x}_i, y_i) are resampled K times, respectively, and the bootstrap samples (\mathbf{x}_i^*, y_i^*) are used to produce estimates $\hat{\boldsymbol{\theta}}^{*(1)}, \dots, \hat{\boldsymbol{\theta}}^{*(K)}$ of $\boldsymbol{\theta}$. The paired bootstrap makes no assumption about error variance homogeneity, and is therefore more robust.

An alternative method, which resamples the estimating equations instead of resampling the data directly, was proposed by Hu and Zidek (1995). For a linear model $\mathbf{Y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{e}$, where $\boldsymbol{\theta} \in \mathbb{R}^p$, $\text{Cov}(\mathbf{e}) = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$ and $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$, the normal estimating equations for the ordinary least-squares estimate $\hat{\boldsymbol{\theta}}$ solve $\sum_{i=1}^n \mathbf{x}_i r_i = 0$, where r_i 's are the residuals. This can be reexpressed as

$$\hat{\boldsymbol{\theta}} = \boldsymbol{\theta} + (\mathbf{X}^T \mathbf{X})^{-1} \sum_{i=1}^n \mathbf{x}_i e_i. \quad (2)$$

The bootstrap version of (2) becomes

$$\hat{\boldsymbol{\theta}}^* = \hat{\boldsymbol{\theta}} + (\mathbf{X}^T \mathbf{X})^{-1} \sum_{i=1}^n \mathbf{z}_i^*, \quad (3)$$

where \mathbf{z}_i^* is a bootstrap sample from $(\mathbf{x}_1 r_1, \dots, \mathbf{x}_n r_n)$. This method is robust against nonhomogeneity of the second and higher moments.

A common feature of all bootstrap methods is that the same estimation procedure is applied to each of the bootstrapped data sets, making the bootstrap computationally expensive, especially for moderate to large p . For example, the number of operations for the paired bootstrap using the standard least-squares algorithm is bounded below by Knp^2 , where p is the number of parameters and K is the number of bootstrap samples.

He and Hu (2002) introduced a new bootstrap method called the Markov chain marginal bootstrap (MCMB), which reduces the computational complexity of the bootstrap by resampling the marginal estimating equations p times at each bootstrap step k to sequentially solve p one-dimensional equations instead of the p -dimensional system of equations. It uses a partially updated vector of parameter estimates

$$\hat{\boldsymbol{\theta}}^{(k,k-1)} = \left(\hat{\theta}_1^{(k)}, \dots, \hat{\theta}_{j-1}^{(k)}, \hat{\theta}_j^{(k)}, \hat{\theta}_{j+1}^{(k-1)}, \dots, \hat{\theta}_p^{(k-1)} \right)$$

to obtain the most current estimate $\hat{\theta}_j^{(k)}$. In Kocherginsky et al. (2005), we applied the MCMB algorithm to quantile regression models (Koenker and Bassett, 1978) and showed that it is a fast and reliable resampling method for problems with large n and p . We also proposed the MCMB-A algorithm, where an affine transformation \mathbf{A} of the parameter space is used to alleviate the problem of possible high autocorrelation of the MCMB sequences. In the MCMB-A algorithm the \mathbf{A} -transformation involves implementing the original MCMB algorithm on a different parameter scale using a transformed parameter vector. While in that article we discussed the application of the \mathbf{A} -transformation in quantile regression, the autocorrelation of MCMB sequences is an even more common and persistent problem in nonlinear regression models.

In the next three sections we review the MCMB algorithm, discuss the MCMB-A method, and describe a further extension of the MCMB algorithm, the \mathbf{B} -transformation, which is a transformation of the estimating equations, aiming to broaden the applicability of the MCMB algorithm to general estimating equations. With both \mathbf{A} and \mathbf{B} -transformations, we show in Section 5 that the MCMB algorithm may be used for inference related to a very general class of estimating equations. We discuss the relationship between the MCMB algorithm and Markov chain Monte Carlo (MCMC) algorithms in Section 6. We then give illustrative examples in Section 7 to show some good properties of the MCMB method, followed by some further remarks in Section 8.

2. Review of MCMB

Let Y_1, \dots, Y_n be a sequence of independent random vectors, and $\theta \in \mathbb{R}^p$ be an unknown p -dimensional parameter. Suppose that $\psi_i(Y_i, \theta)$ are \mathbb{R}^p -valued functions of the unknown parameter $\theta \in \mathbb{R}^p$ and the observations Y_i . Assume that the expected value of $\psi_i(Y_i, \theta)$'s is the p -dimensional vector of 0's, i.e., $E_{\theta}[\psi_i(Y_i, \theta)] = \mathbf{0}$. Consider $\hat{\theta}$, a solution of the unbiased estimation equation

$$S(\mathbf{Y}, \theta) = \sum_{i=1}^n \psi_i(Y_i, \theta) = \mathbf{0}, \quad (4)$$

where \mathbf{Y} denotes the data, and the subscript n on the estimator $\hat{\theta}$ is suppressed. Using this notation, $S(\mathbf{Y}, \theta)$ is a p -dimensional vector. Its j th element is $S_j(\mathbf{Y}, \theta) = \sum_{i=1}^n \psi_{ij}(Y_i, \theta)$, $j = 1, \dots, p$, and $\psi_i(Y_i, \theta)$ is a vector $(\psi_{i1}(Y_i, \theta), \dots, \psi_{ip}(Y_i, \theta))^T$. Assume that a decomposition $\psi_i(Y_i, \theta) = \mathbf{a}_i z_i$ exists, so that \mathbf{a}_i 's are vectors or matrices independent of \mathbf{Y} , and z_i 's are random variables depending on \mathbf{Y} and θ . For example, in a linear regression model $Y_i = \mathbf{x}_i^T \theta + e_i$, the MCMB can be carried out with $\mathbf{a}_i = \mathbf{x}_i$ and $z_i = e_i$. The \mathbf{a}_i 's and z_i 's are used as a matter of convenient notation to indicate the quantities that are to be resampled; they do not provide a unique decomposition, allowing for different resampling schemes. Note that $\hat{\theta}$ is an M -estimator. Typically, M -estimators (Huber, 1964) minimize $\sum_{i=1}^n \rho(y_i, \theta)$ for an objective function ρ , often chosen to be a symmetric and continuous differentiable function. If ρ is differentiable with $\rho' = \psi$, then the M -estimator is defined as the solution of $\sum_{i=1}^n \psi(y_i, \theta) = \mathbf{0}$.

The conventional bootstrap solves (4) for each bootstrap sample, which is obtained with residual or paired bootstrap, or by sampling directly from the estimating equations. Solving for $\hat{\theta}$ involves solving the p -dimensional system of Eq. (4) a large number of times. The MCMB instead solves the sequences of one-dimensional equations using the following algorithm:

1. Initialize $\hat{\theta}^{(0)} = \hat{\theta}$ and $k = 1$.
2. (a) Set $j = 1$.
(b) Draw a bootstrap sample $\{z_{1j}^{*(k)}, \dots, z_{nj}^{*(k)}\}$ from $\{z_1, \dots, z_n\}$ for each $j = 1, \dots, p$.
3. (a) Solve for $\hat{\theta}_j^{(k)}$

$$S_j\left(\mathbf{Y}; \hat{\theta}_1^{(k)}, \dots, \hat{\theta}_{j-1}^{(k)}, \boxed{\hat{\theta}_j^{(k)}}, \hat{\theta}_{j+1}^{(k-1)}, \dots, \hat{\theta}_p^{(k-1)}\right) = S_j^{*(k)}\left(\mathbf{Y}; \hat{\theta}^{(0)}\right), \quad (5)$$

where $(\hat{\theta}_1^{(k)}, \dots, \hat{\theta}_{j-1}^{(k)}, \hat{\theta}_{j+1}^{(k-1)}, \dots, \hat{\theta}_p^{(k-1)})$ is the most current estimate of $\hat{\theta}$ without the j th component, and

$$S_j^{*(k)}(\mathbf{Y}; \hat{\theta}^{(0)}) = \sum_{i=1}^n \mathbf{a}_i z_{ij}^{*(k)}.$$

- (b) Return to 2(a) with $j = j + 1$, if $j < p$.
4. Increase k by 1 and go to step 2a or stop if k has reached a pre-specified level K .

This algorithm results in a sequence of solutions $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(K)}$, which by construction is a Markov chain. At the k th iteration the MCMB algorithm solves p one-dimensional equations for $\hat{\theta}_j^{(k)}$, $j = 1, \dots, p$ using the already updated components $\hat{\theta}_1^{(k)}, \dots, \hat{\theta}_{j-1}^{(k)}$ together with the values remaining from the previous iteration $\hat{\theta}_{j+1}^{(k-1)}, \dots, \hat{\theta}_p^{(k-1)}$.

3. Discussion on MCMB-A

Recall that MCMB solves the unbiased estimation equations (4) for θ . If the covariance matrix of $S(\mathbf{Y}, \theta)$ is nearly singular, the resulting MCMB sequences will be strongly autocorrelated. We propose a transformation of the parameter space, which we call the **A-transformation**, which “decorrelates” the MCMB sequences. Accordingly, consider a $p \times p$ matrix \mathbf{A} and a vector $\tilde{\theta}$ such that

$$\theta = \mathbf{A}\tilde{\theta},$$

and consequently $\tilde{\theta} = \mathbf{A}^{-1}\theta$. The transformed estimating equations are

$$\tilde{S}(\mathbf{Y}, \tilde{\theta}) = \sum_{i=1}^n \frac{\partial}{\partial \tilde{\theta}} \rho(y_i, \mathbf{A}\tilde{\theta}) = \sum_{i=1}^n \frac{\partial \theta}{\partial \tilde{\theta}} \frac{\partial}{\partial \theta} \rho(y_i, \mathbf{A}\tilde{\theta}) = \mathbf{A} \cdot S(\mathbf{Y}, \theta).$$

In order to “decorrelate” the sequences of MCMB solutions, the covariance matrix of $\tilde{S}(\mathbf{Y}, \tilde{\theta})$ should be aimed at the identity matrix:

$$\text{Cov}(\tilde{S}(\mathbf{Y}, \tilde{\theta})) = \mathbf{A} \text{Cov}(S(\mathbf{Y}, \theta)) \mathbf{A}^T = \mathbf{I}.$$

Solving this equation results in

$$\mathbf{A} = \text{Cov}(S(\mathbf{Y}, \theta))^{-1/2}. \quad (6)$$

Note that \mathbf{A} exists since $\text{Cov}(S(\mathbf{Y}, \theta))$ is symmetric and is assumed to be invertible, and its square root can be found using the singular value decomposition. The covariance matrix $\text{Cov}(S(\mathbf{Y}, \theta))$ can be estimated by

$$\text{Cov}(S(\mathbf{Y}, \theta)) \approx \frac{1}{n} \{\psi(\mathbf{Y}, \hat{\theta})^T \psi(\mathbf{Y}, \hat{\theta})\},$$

where $\hat{\theta}$ is the initial estimate of θ , and $\psi(\mathbf{Y}, \hat{\theta})$ is an $n \times p$ matrix with elements $\psi_{ij}(y_i, \hat{\theta})$.

Let $S(\mathbf{Y}, \theta)$ be such that $\partial/\partial \theta S(\mathbf{Y}, \theta)$ is continuous in θ , and the matrix $E(\partial/\partial \theta S(\mathbf{Y}, \theta))$ is nonsingular. If $\hat{\theta}$ is a consistent estimator of θ and solves $S(\mathbf{Y}, \theta) = \mathbf{0}$, then, from the Taylor expansion,

$$\mathbf{A}\text{-Var}(\hat{\theta}) = \left(E \frac{\partial}{\partial \theta} S(\mathbf{Y}, \theta) \right)^{-1} \text{Var}(S(\mathbf{Y}, \theta)) \left(E \frac{\partial}{\partial \theta} S(\mathbf{Y}, \theta) \right)^{-1},$$

where $\mathbf{A}\text{-Var}(\hat{\theta})$ is the asymptotic variance of $\hat{\theta}$. Furthermore, if $\text{Var}(S(\mathbf{Y}, \theta)) = \pm E(\partial/\partial \theta S(\mathbf{Y}, \theta))$, then $\mathbf{A}\text{-Var}(\hat{\theta}) = \text{Var}(S(\mathbf{Y}, \theta))^{-1}$. Also, it follows from (6) that

$$\begin{aligned} \mathbf{A}\text{-Var}(\hat{\theta}) &= \mathbf{A}^{-1} \left(\frac{\partial}{\partial \theta} S(\mathbf{Y}, \theta) \right)^{-1} \text{Var}(S(\mathbf{Y}, \theta)) \left(\frac{\partial}{\partial \theta} S(\mathbf{Y}, \theta) \right)^{-1} \mathbf{A}^{-1} \\ &= \text{Var}(S(\mathbf{Y}, \theta))^{1/2} \text{Var}(S(\mathbf{Y}, \theta))^{-1} \text{Var}(S(\mathbf{Y}, \theta))^{1/2} = \mathbf{I} \end{aligned}$$

when $E(\partial/\partial \theta S(\mathbf{Y}, \theta)) = \text{Var}(S(\mathbf{Y}, \theta))^{-1}$. Thus, the $\tilde{\theta}_j$'s are uncorrelated, and the variance of $\hat{\theta}_j$ is 1. The latter is a useful property, since it provides the scale for the variance of $\tilde{\theta}$ when the initial estimate may not be accurate.

In practice, any consistent estimate of \mathbf{A} up to a scale factor may be used for the same purpose of decorrelating the MCMB sequence. For this reason, the \mathbf{A} -transformation for a linear model can be used in the following simple way:

$$\mathbf{X}\theta = \mathbf{X}\mathbf{A}\tilde{\theta} = \tilde{\mathbf{X}}\tilde{\theta},$$

where $\tilde{\mathbf{X}} = \mathbf{X}\mathbf{A}$, and $r_i = y_i - \mathbf{X}\theta = y_i - \tilde{\mathbf{X}}\tilde{\theta}$. In this special case, the MCMB algorithm can be applied directly on $\tilde{\theta}$ scale if the design matrix \mathbf{X} is replaced by $\tilde{\mathbf{X}}$ in the beginning, and transformed back in the end. The transformation between $\tilde{\theta}$ and θ becomes unnecessary at every step.

4. Introduction to MCMB-B

He and Hu (2002) provided the conditions under which the distribution of $\hat{\theta}^{(k)}$ resulting from the MCMB algorithm is an asymptotically valid approximation to that of $\hat{\theta}$. Again, we consider an estimator defined by (4). The derivations in He and Hu (2002) show that if the following conditions hold, the MCMB sequence does approximate the distribution of $\hat{\theta}$ as desired.

(C1) $\partial \psi_i(y_i, \theta)/\partial \theta$ is continuously differentiable in θ ;

(C2) $n^{-1} \sum_{i=1}^n (\partial^2/\partial \theta^2) \psi_i(Y_i, \theta)$ is uniformly bounded in n , and $\sum_{i=1}^n E|\psi_i(Y_i, \theta)|^3 = \mathcal{O}(n)$ as $n \rightarrow \infty$;

(C3) both $n^{-1} \sum_{i=1}^n [\psi_i(Y_i, \theta) \psi_i(Y_i, \theta)^T]$ and $-n^{-1} \sum_{i=1}^n (\partial/\partial \theta) \psi_i(Y_i, \theta)$ converge to the information matrix $D(\theta) > 0$ as $n \rightarrow \infty$.

Condition (C3) is essential for $\sqrt{n}(\hat{\theta} - \hat{\theta}^{(k)})$ to be an approximate multiple AR(1) process. This condition holds if (4) is a likelihood equation and Y_i are i.i.d. with $\psi_i = \psi$ for all i , where D is the Fisher information matrix. In general, however we cannot expect (C3) to hold.

We propose a transformation of the estimating equations which ensures that (C3) holds. We call it the **B**-transformation, as it “broadens” the applicability of the MCMB. For any estimator $\hat{\theta}$ which is the solution of $\psi(\mathbf{Y}, \theta) = S(Y, \theta) = \sum_{i=1}^n \psi_i(Y_i, \theta) = \mathbf{0}$, where $\psi(\mathbf{Y}, \theta)$ is a p -dimensional vector consider a linear transformation of the estimating equations

$$\tilde{\psi}(\mathbf{Y}, \theta) = \mathbf{B}\psi(\mathbf{Y}, \theta),$$

where \mathbf{B} is a $p \times p$ nonsingular matrix and $\tilde{\psi}(\mathbf{Y}, \theta)$ is a $p \times 1$ vector. Multiplying the original set of equations by a nonsingular matrix transforms the marginal equations, but does not affect the solution $\hat{\theta}$. The new equations have mean and variance

$$E \frac{\partial}{\partial \theta} \tilde{\psi}(\mathbf{Y}, \theta) = \mathbf{B} E \frac{\partial}{\partial \theta} \psi(\mathbf{Y}, \theta), \quad \text{Var}(\tilde{\psi}(\mathbf{Y}, \theta)) = \mathbf{B} \text{Var}(\psi(\mathbf{Y}, \theta)) \mathbf{B}^T. \quad (7)$$

In order for the condition (C3) to hold under the transformed equations, we require

$$-E_{\theta} \frac{\partial \tilde{\psi}(\mathbf{Y}, \theta)}{\partial \theta} = \text{Var}(\tilde{\psi}(\mathbf{Y}, \theta)). \quad (8)$$

Substituting (7) into (8) and solving for \mathbf{B} results in

$$\mathbf{B} = -E_{\theta} \left[\frac{\partial \psi(\mathbf{Y}, \theta)}{\partial \theta} \right]^T \cdot [\text{Var}(\psi(\mathbf{Y}, \theta))]^{-1},$$

which ensures that (C3) holds for the new equations. Note that if (C3) holds for the original equations $\psi(\mathbf{Y}, \theta)$ then $\mathbf{B} = \mathbf{I}$ and $\tilde{\psi}(\mathbf{Y}, \theta) = \psi(\mathbf{Y}, \theta)$.

The quantities $E_{\theta}[\partial/\partial \theta \psi(\mathbf{Y}, \theta)]$ and $\text{Var}(\psi(\mathbf{Y}, \theta))$ can be obtained using the initial estimate $\hat{\theta}$ of θ . Also note that (C3) needs to hold only for the true parameter θ . If \mathbf{B} is evaluated at $\hat{\theta}$, which is a consistent estimate of θ , (8) holds asymptotically. We refer to this estimation procedure as MCMB-**B**.

5. Journey to MCMB-AB

Although the **B**-transformation enforces (C3), it does nothing to eliminate the autocorrelation of the resulting MCMB sequences. We therefore propose to use both **A**- and **B**-transformations jointly to address both issues simultaneously, and call the resulting algorithm MCMB-**AB**.

Let matrices \mathbf{P} and \mathbf{Q} be

$$\mathbf{P} = E \left[\frac{\partial}{\partial \theta} \psi(\mathbf{Y}, \theta) \right]^T \quad \text{and} \quad \mathbf{Q} = E[\psi(\mathbf{Y}, \theta)\psi(\mathbf{Y}, \theta)^T]. \quad (9)$$

The possible dependence of \mathbf{P} and \mathbf{Q} on n will be suppressed. Also, let $\theta = \mathbf{A}\tilde{\theta}$ and $\tilde{\psi}(\mathbf{Y}, \tilde{\theta}) = \mathbf{A}\psi(\mathbf{Y}, \mathbf{A}\tilde{\theta}) = \mathbf{0}$, as before. The new estimating equations, which result from applying the **B**-transformation after the **A**-transformation, become

$$\mathbf{B}\mathbf{A}\psi(\mathbf{Y}, \mathbf{A}\tilde{\theta}) = \mathbf{0}. \quad (10)$$

We choose \mathbf{B} so that condition (C3) holds for the new equations in (10), that is

$$E \left[\frac{\partial}{\partial \tilde{\theta}} (\mathbf{B}\mathbf{A}\psi(\mathbf{Y}, \mathbf{A}\tilde{\theta})) \right] = -E[(\mathbf{B}\mathbf{A}\psi(\mathbf{Y}, \theta)) \cdot (\mathbf{B}\mathbf{A}\psi(\mathbf{Y}, \theta))^T]. \quad (11)$$

The expression on the left side of (11) simplifies to $\mathbf{B} \mathbf{A} E[(\partial/\partial \theta) \psi(\mathbf{Y}, \theta)] \cdot \partial \theta / \partial \tilde{\theta} = \mathbf{B} \mathbf{A} \mathbf{P} \mathbf{A}$, and on the right side simplifies to $-\mathbf{B} \mathbf{A} \mathbf{Q} \mathbf{A}^T \mathbf{B}^T$. Equating the two and solving for \mathbf{B} results in

$$\mathbf{B} = -\mathbf{A} \mathbf{P} \mathbf{Q}^{-1} \mathbf{A}^{-1}. \quad (12)$$

Next, we choose \mathbf{A} to eliminate the autocorrelation, so that

$$\text{Var}(\mathbf{B}\mathbf{A}\psi(\mathbf{Y}, \theta)) = \mathbf{I}.$$

Substituting \mathbf{B} from (12), and solving for \mathbf{A} results in

$$\mathbf{A} = (n \cdot \mathbf{P}\mathbf{Q}^{-1}\mathbf{P})^{-1/2}. \quad (13)$$

Note that if $\tilde{\theta}$ is the solution of $\mathbf{B}\tilde{\psi}(\mathbf{Y}, \tilde{\theta}) = \mathbf{0}$, then the asymptotic covariance is $\mathbf{A}\text{-Var}(\hat{\theta}) = \mathbf{I}$, since

$$\begin{aligned} \mathbf{A}\text{-Var}(\hat{\theta}) &= E\left(\frac{\partial}{\partial \theta} \mathbf{B}\mathbf{A}S(\mathbf{Y}, \theta)\right)^{-1} \cdot \text{Var}(\mathbf{B}\mathbf{A}S(\mathbf{Y}, \theta)) \cdot E\left(\frac{\partial}{\partial \theta} \mathbf{B}\mathbf{A}S(\mathbf{Y}, \theta)\right)^{-1} \\ &= \text{Var}(\mathbf{B}\mathbf{A}S(\mathbf{Y}, \theta))^{-1} = \mathbf{I}, \end{aligned}$$

where the second equality holds because \mathbf{B} is chosen so that $\text{Var}(\mathbf{B}\mathbf{A}S(\mathbf{Y}, \theta)) = -E(\partial/\partial \theta \mathbf{B}\mathbf{A}S(\mathbf{Y}, \theta))$.

In practice, since the true parameter θ is unknown, \mathbf{A} and \mathbf{B} are estimated using $\hat{\theta}$, a consistent estimate of θ . Then, the $p \times p$ matrices \mathbf{P} and \mathbf{Q} are computed according to

$$\hat{\mathbf{P}} = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \psi_i(Y_i, \theta) \bigg|_{\theta=\hat{\theta}} \quad \text{and} \quad \hat{\mathbf{Q}} = \frac{1}{n} \sum_{i=1}^n \left\{ \psi_i(Y_i, \hat{\theta}) \psi_i(Y_i, \hat{\theta})^T \right\}. \quad (14)$$

The matrices \mathbf{A} and \mathbf{B} are then computed according to (13) and (12) with \mathbf{P} and \mathbf{Q} replaced by $\hat{\mathbf{P}}$ and $\hat{\mathbf{Q}}$, that is $\hat{\mathbf{A}} = (n \cdot \hat{\mathbf{P}}\hat{\mathbf{Q}}^{-1}\hat{\mathbf{P}})^{-1/2}$ and $\hat{\mathbf{B}} = -\hat{\mathbf{A}}\hat{\mathbf{P}}\hat{\mathbf{Q}}^{-1}\hat{\mathbf{A}}^{-1}$.

The following proposition shows the asymptotic validity of MCMB-AB. We follow the convention to say that $\{\psi_{ij}(\cdot, \theta)\}$ is continuous in θ uniformly in i, j if for any $\varepsilon > 0$ there exists $\delta > 0$ such that $\max_{i,j} |\psi_{ij}(\cdot, \theta_1) - \psi_{ij}(\cdot, \theta_2)| < \varepsilon$ whenever $\|\theta_1 - \theta_2\| < \delta$.

Proposition 1. *If $\psi_{ij}(Y_i, \theta)$ and $\partial \psi_{ij}(Y_i, \theta)/\partial \theta$ are continuous in θ uniformly in i, j , and if the matrices \mathbf{P} and \mathbf{Q} of (9) exist and have finite limits \mathbf{P}_∞ and \mathbf{Q}_∞ as $n \rightarrow \infty$, then condition (C3) is satisfied when $\psi_i(Y_i, \theta)$ is replaced by $\hat{\mathbf{B}}\hat{\mathbf{A}}\psi_i(Y_i, \theta)$.*

Proof. Let θ_0 denote the true parameter. For convenience, let $\mathbf{A}^* = (\mathbf{P}_\infty \mathbf{Q}_\infty^{-1} \mathbf{P}_\infty)^{-1/2}$ and $\mathbf{B}^* = -\mathbf{A}^* \mathbf{P}_\infty \mathbf{Q}_\infty^{-1} \mathbf{A}^{*-1}$. Since $\hat{\theta}$ is a consistent estimator of θ_0 , and by continuity of $\psi(\mathbf{Y}, \theta)$ and its derivative with respect to θ , we have

$$\begin{aligned} \hat{\mathbf{P}} - \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \psi_i(Y_i, \theta) \bigg|_{\theta=\theta_0} &\rightarrow 0, \\ \hat{\mathbf{Q}} - \frac{1}{n} \sum_{i=1}^n \psi_i(Y_i, \theta_0) \psi_i(Y_i, \theta_0)^T &\rightarrow 0. \end{aligned} \quad (15)$$

By the law of large numbers,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \psi_i(Y_i, \theta) \bigg|_{\theta=\theta_0} - \mathbf{P} &\rightarrow 0, \\ \frac{1}{n} \sum_{i=1}^n \psi_i(Y_i, \theta_0) \psi_i(Y_i, \theta_0)^T - \mathbf{Q} &\rightarrow 0, \end{aligned}$$

and thus $\hat{\mathbf{P}} - \mathbf{P}_\infty \rightarrow 0$ and $\hat{\mathbf{Q}} - \mathbf{Q}_\infty \rightarrow 0$. Therefore, $\sqrt{n} \hat{\mathbf{A}} \rightarrow \mathbf{A}^*$ and $\hat{\mathbf{B}} \rightarrow \mathbf{B}^*$ as $n \rightarrow \infty$. By the construction of \mathbf{A}^* and \mathbf{B}^* ,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (\mathbf{B}^* \mathbf{A}^* \psi_i(Y_i, \theta_0)) (\mathbf{B}^* \mathbf{A}^* \psi_i(Y_i, \theta_0))^T &\rightarrow \mathbf{B}^* \mathbf{A}^* \mathbf{Q}_\infty \mathbf{A}^{*T} \mathbf{B}^{*T}, \\ -\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \mathbf{B}^* \mathbf{A}^* \psi_i(Y_i, \theta) \bigg|_{\theta=\theta_0} &\rightarrow -\mathbf{B}^* \mathbf{A}^* \mathbf{P}_\infty. \end{aligned}$$

These limits also hold with $(\mathbf{A}^* \mathbf{B}^*)$ replaced by $\sqrt{(n)} \hat{\mathbf{A}} \hat{\mathbf{B}}$, and $\mathbf{D}(\theta_0) = \mathbf{B}^* \mathbf{A}^* \mathbf{Q}_\infty \mathbf{A}^{*T} \mathbf{B}^{*T}$. \square

Note that uniform continuity assumption in the proposition is only used to simplify the proof that leads to (15). In many cases, the results in (15) hold without such strong conditions. Under mild technical conditions, the MCMB-**AB** method remains valid if ψ has finitely many jumps, as in the quantile regression estimators. In such cases, the derivative of ψ in (14) needs to be replaced by derivative of the expected value of ψ .

The MCMB-**AB** method extends the MCMB approach to general estimating equations by exploring the flexibility of redefining margins in a class of equivalent p -dimensional equations. In some cases, this would result in more complicated marginal equations, but one-dimensional equations are generally manageable even if a grid search algorithm is necessary to find a reliable root.

6. An MCMC perspective

The MCMB method shares one important feature with the Gibbs sampler, a popular MCMC algorithm, that is, both methods break up a high-dimensional problem into one-dimensional problems. In the asymptotic sense, one may view the update on the j th component of θ as an attempt to draw from the conditional distribution of θ_j given the current values of θ without the j th component, but the exact asymptotic equivalence here is difficult to specify. The differences between the MCMB method and the MCMC method are clearer, because the latter is a finite-sample procedure and the former is based on asymptotic normality of the estimators. As a consequence, only a short chain from the MCMB algorithm is usually needed to achieve convergence while a long chain is preferred for any MCMC method. This is a double-edged sword. Convergence of MCMC algorithms is often a source of complaint, but accuracy from any MCMC method generally increases with chain length. On the other hand, the MCMB algorithm starts with a consistent estimator of θ and converges rather quickly, but error accumulation can occur when a long chain is run. In fact, we advice against using a long MCMB chain, because the accuracy of the MCMB algorithm is usually dominated by the normal approximation to the distribution of $\hat{\theta}$.

A recent paper of Tian et al. (2007) provides a new perspective on the relationship between the MCMB method, the estimating bootstrap method (Hu and Zidek, 1995), and the MCMC approach to estimating equations. Tian et al. (2007) uses an MCMC algorithm to the asymptotic distribution of $\hat{\theta}$ based on a different **B**-transformation (from what is used in this paper) to achieve asymptotic normality of the estimating equation at the true parameter.

7. Illustrative examples

7.1. A simulation study

To illustrate the use of **A**- and **B**-transformations, we consider a simple two-parameter nonlinear regression model with a heteroscedastic error term

$$y_i = \theta_1 e^{-\theta_2 x_i} + e^{x_i} z_i \quad i = 1, \dots, 500, \quad (16)$$

where $-1 < x_i < 2$, and $z_i \sim N(0, 1)$. The scatter plot of a typical data set generated from this model with $\theta_1 = 4$ and $\theta_2 = 1.5$ can be seen in Fig. 1. For this model **P** and **Q** are

$$\mathbf{P} = E \begin{pmatrix} -e^{-2x\theta_2} & x e^{-\theta_2 x} (\theta_1 e^{-\theta_2 x} - z) \\ x e^{-\theta_2 x} (\theta_1 e^{-\theta_2 x} - z) & -\theta_1 x^2 e^{-\theta_2 x} (\theta_1 e^{-\theta_2 x} - z) \end{pmatrix},$$

$$\mathbf{Q} = E \begin{pmatrix} z^2 e^{-2\theta_2 x} & -\theta_1 x z^2 e^{-2\theta_2 x} \\ -\theta_1 x z^2 e^{-2\theta_2 x} & \theta_1^2 x^2 z^2 e^{-2\theta_2 x} \end{pmatrix}.$$

We can empirically check whether condition (C3) holds by checking whether $\hat{\mathbf{P}}\hat{\mathbf{Q}}^{-1} \approx \mathbf{I}$. For the data set in Fig. 1 we find that

$$\hat{\mathbf{P}}\hat{\mathbf{Q}}^{-1} = \begin{pmatrix} -1.823 & -0.542 \\ -4.140 & -1.904 \end{pmatrix}$$

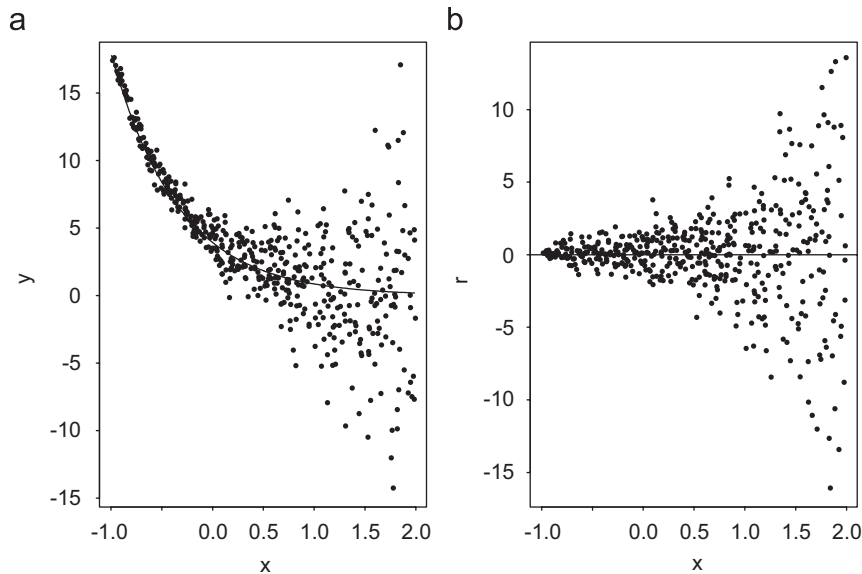


Fig. 1. (a) Scatter plot of the data in (16), (b) residuals from `nls()` vs. x .

with eigenvalues $\lambda_1 = -3.36, \lambda_2 = -0.37$. Recall that if the matrix is approximately equal to the identity matrix, the ratio of the largest and the smallest eigenvalues (also called the condition number of a matrix) should be close to 1. The condition number of in $\hat{\mathbf{P}}\hat{\mathbf{Q}}^{-1}$ is 9.22, suggesting that (C3) does not hold.

We compared the performance of the following methods: nonlinear least-squares regression (Bates and Watts, 1988) using `nls()` available in **R** (NLS), direct MCMB (MCMB), MCMB with the **A**-transformation only (MCMB-A), MCMB with the **B**-transformation only (MCMB-B), MCMB with the **A**- and **B**-transformations (MCMB-AB), the paired bootstrap (PB), and the Monte Carlo simulation (MC). All bootstrap estimates were based on $K = 500$ replications. In the Monte Carlo simulation $M = 500$ independent \mathbf{y} 's were generated, and `nls()` was used to estimate the parameters.

The Markov chains $\hat{\theta}^{(k)}$ produced by each of the four MCMB methods are plotted in Fig. 2. Applying MCMB directly generates Markov chains with some visible autocorrelation (a), suggesting that the **A**-transformation is needed. Panel (b) shows Markov chains produced MCMB with only the **A**-transformation. Although the resulting sequences look somewhat improved, clearly autocorrelation persists. Panel (c) shows MCMB chains generated by MCMB-B, which clearly performs substantially worse than direct MCMB or MCMB-A. Based on this and other models (not shown), we do not recommend using the **B**-transformation alone. Finally, panel (d) shows Markov chains resulting from MCMB-AB, where the sequences look completely random.

Mean and standard error estimates are reported in Table 1. Table 2 shows the corresponding coverage probabilities based on the asymptotic $(1 - \alpha)100\%$ confidence intervals. Note that the estimated MCMB and MCMB-A standard errors are substantially smaller than in MC or PB methods, resulting in much lower than nominal coverage probabilities. However, using both transformations jointly dramatically improves the MCMB performance. MCMB-AB performs as well as the paired bootstrap, providing accurate coverage probabilities and parameter estimates close to the true θ . We also note the poor performance of the usual NLS estimators, which can be improved upon with either an appropriate transformation of the response or by using the robust sandwich variance-covariance estimator.

Fig. 2 and Table 1 both demonstrate that MCMB-AB improves the MCMB performance when both autocorrelation and heteroscedasticity are present. Clearly, MCMB-B had the worst performance of the four MCMB methods, and can introduce even worse autocorrelation in the MCMB chains than direct MCMB.

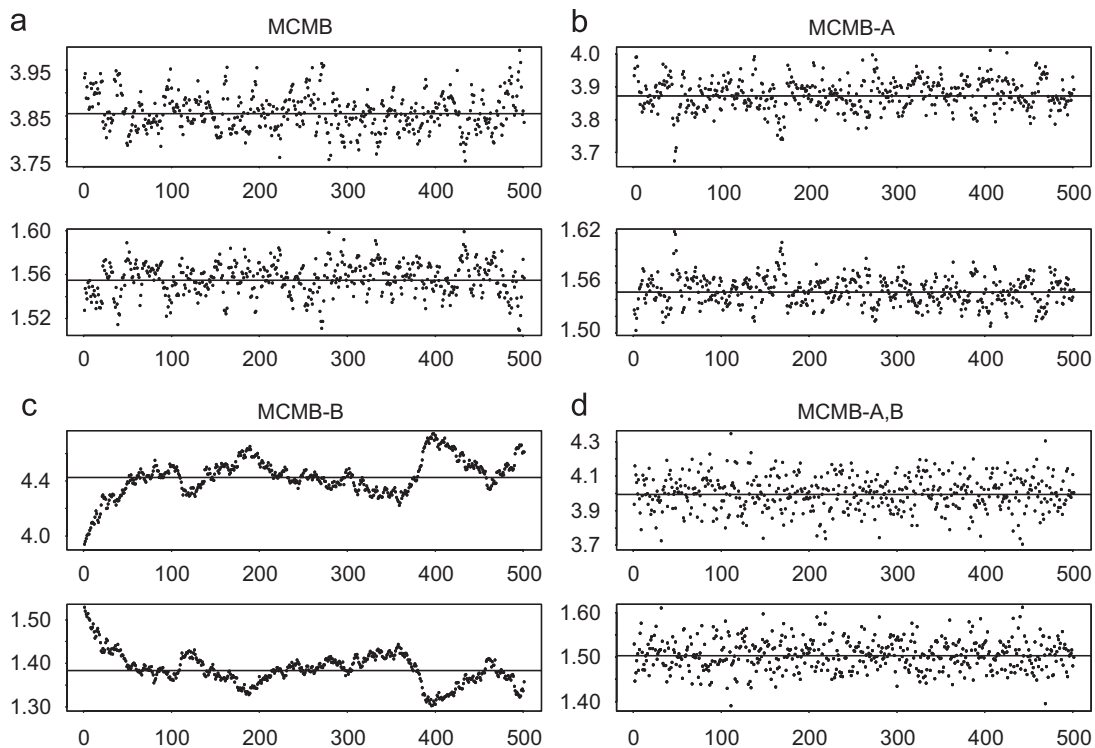


Fig. 2. (a) MCMB chains, (b) MCMB-A, (c) MCMB-B, (d) MCMB-AB. In each subfigure, the upper panel is for θ_1 and the lower panel is for θ_2 .

Table 1
Parameter estimates and their standard errors

	$\hat{\theta}_1$	$SE(\hat{\theta}_1)$	$\hat{\theta}_2$	$SE(\hat{\theta}_2)$
NLS	3.8698	0.2017	1.5497	0.0726
MCMB	3.8551	0.0404	1.5549	0.0148
MCMB-A	3.8724	0.0475	1.5488	0.0171
MCMB-B	4.4246	0.1329	1.3838	0.0366
MCMB-AB	3.9943	0.0994	1.5029	0.0345
PB	3.8790	0.0984	1.5470	0.0346
MC	3.9978	0.1058	1.5001	0.0370

We also note that the confidence intervals can often be improved upon by using the bias-corrected and accelerated (BC_a) version of the percentile method for confidence interval estimation (Efron and Tibshirani, 1998), but it is yet to be determined whether BC_a will benefit the MCMB methods.

7.2. An example

We now demonstrate another property of the MCMB method in a relatively simple case: the variance of median regression estimates with leverage points in the data.

Suppose that the median of y given x is $\alpha + \beta x$, which is estimated by minimizing $\sum_{i=1}^n |y_i - \alpha - \beta x_i|$ for a given sample $\{(x_i, y_i)\} = \{(0, 1), (1, 1), (1, 2), (2, 3), (3, 2), (10, 10)\}$. This artificial data set includes a leverage point at $(10, 10)$. The usual paired bootstrap based on 100 bootstrap samples gave the standard error estimate of 0.36 for the estimate of β , while the MCMB method based on a chain of length 100 estimated the standard error to be 0.08. Both methods avoid the specification of likelihood.

Table 2
Coverage probability and average confidence interval length

	$\alpha = 0.1$				$\alpha = 0.05$			
	Coverage		Length		Coverage		Length	
	θ_1	θ_2	θ_1	θ_2	θ_1	θ_2	θ_1	θ_2
MCMB	0.38	0.41	0.3121	0.1094	0.45	0.50	0.3719	0.1303
MCMB-A	0.47	0.48	0.3691	0.1272	0.57	0.54	0.4398	0.1515
MCMB-AB	0.89	0.89	0.8734	0.3134	0.93	0.95	1.0407	0.3734
PB	0.90	0.92	0.9075	0.3177	0.94	0.94	1.0814	0.3786
NLS	0.97	0.96	1.4862	0.5214	1.00	1.00	1.7710	0.6212

The significant difference in the results here can be easily attributed to the leverage point in the sample, but if the linear quantile model is true and the conditional distributions of y given x do not change much with x , the paired bootstrap inflates the variance. This problem with the usual bootstrap has been observed even in the univariate sample median studies, see for example, [Ghosh et al. \(1984\)](#) and [Jimenez-Gamero et al. \(1998\)](#). The MCMB method uses all the design points, so it captures the conditional regression analysis in a natural way. For algorithmic development of the MCMB method in quantile regression, we refer to [Kocherginsky et al. \(2005\)](#).

8. Conclusions

In this article we introduced a transformation of the estimating equations, called the **B**-transformation, which validates the MCMB for estimating equations that are not necessarily likelihood-based. We recommend that this transformation be used jointly with the **A**-transformation of the parameter space, which “decorrelates” the MCMB sequences.

The MCMB-AB approach for inference may be used as an alternative to the large-sample inference based on the robust variance–covariance matrix for any estimator given by a generalized estimating equation. The purpose of the MCMB approach (including the extensions discussed here) is not to improve accuracy over the usual bootstrap method. Instead, we focus on two distinctive features of the MCMB approach.

First, when the estimating equation is difficult to solve (especially when p is modestly large), the usual bootstrap is practically prohibitive. For example, some estimating equations are likely to have multiple roots, so finding the right root as a consistent estimator is challenging even in the original sample. In such cases, the MCMB method can be used by relying on a consistent estimator of θ which is easier to compute but is not the same as the estimator defined by the estimating equation.

Second, the MCMB algorithm requires solving one-dimensional equations only, which is feasible even if the equation is highly nonlinear. The grid search might be the last resort. When multiple roots are present, we can always use the one that is closest to a consistent estimator.

However, consistent estimation of the **A**- and **B**-transformations may, in some cases, require as much as what is needed to estimate the robust variance–covariance matrix itself. The relative merits on the MCMB-AB method need to be established case by case. In the cases where **P** involves a nonparametric density function, as in the quantile regression estimation, the MCMB method has been shown to be robust even when the estimate of **P** itself and the estimate of the robust variance–covariance matrix are sensitive to the smoothing parameter needed for the density estimation. More details on the **A**-, **B**-transformations may be found at [Kocherginsky \(2003\)](#), the first author’s dissertation at the University of Illinois.

References

- Bates, D.M., Watts, D.G., 1988. *Nonlinear Regression Analysis and its Applications*. Wiley, New York, NY.
- Davison, A.C., Hinkley, D.V., 1997. *Bootstrap Methods and their Application*. Cambridge University Press, Cambridge.
- Efron, B., Tibshirani, R.J., 1998. *An Introduction to the Bootstrap*. CRC Press LLC, Boca Raton, FL.

- Ghosh, M., Parr, M.C., Singh, K., Babu, G.J., 1984. A note on bootstrapping the sample median. *Ann. Statist.* 12, 1130–1135.
- He, X., Hu, F., 2002. Markov chain marginal bootstrap. *J. Amer. Statist. Assoc.* 97 (459), 783–795.
- Hu, F., Zidek, J.V., 1995. A Bootstrap based on the estimating equations of the linear model. *Biometrika* 82, 263–275.
- Huber, P.J., 1964. Robust estimation of a location parameter. *Ann. Math. Statist.* 35, 73–101.
- Jimenez-Gamero, M.D., Munoz-Garcia, J., Pino-Mejias, R., 1998. Reduced bootstrap for the median. *Statist. Sinica* 14, 1179–1198.
- Kocherginsky, M., 2003. Extensions of Markov chain marginal bootstrap. Ph.D. Thesis, University of Illinois at Urbana-Champaign.
- Kocherginsky, M., He, X., Mu, Y., 2005. Practical confidence intervals for regression quantiles. *J. Comput. Graphical Statist.* 14, 41–55.
- Koenker, R.W., Bassett, G.J., 1978. Regression quantiles. *Econometrica* 46, 33–50.
- Tian, L., Liu, J.S., Wei, L.J., 2007. Implementation of estimating-equation based inference procedures with MCMC samplers. *J. Amer. Statist. Assoc.*, in press.